

# Linux File Systems

Adrien 'schischi' Schildknecht

July 18, 2014

- Overhead: non-volatile memory is often the slowest component
- Reliability: no fault tolerated (data loss!)
- People expect more and more features...
  - Quotas
  - Snapshot
  - Versioning
  - Replication
  - Database features...

# Section 1

## Hardware

Linux File  
Systems

Adrien  
'schischi'  
Schildknecht

- ⊕ no moving heads nor rotating platters
- ⊕ good random access
- ⊕ low power
- ⊖ wear out (writing)
- ⊖ write in 2 phases
  - Clear group of pages, 128Ko+ (slow)
  - Write individual pages, 4Ko (fast)

## HDD:

- locality
- optimize seeking

## SSD:

- TRIM

## Common:

- Smallest unit is a block (reading/writing)
- Least possible writing

## Section 2

# Abstraction

**File:** store a named piece of data to later retrieve it.

- stream of bytes, let the programmer read raw bytes
- structural metadata: inode
- descriptive metadata: attributes (name, owner, ...)

**Directory:** provide a way to organize multiple files

- hierarchies: directories can contain directories
- data structure which contains names and handles



- **FS**: machinery to store/retrieve data
- **block**: smallest unit writable by a disk or fs
- **metadata**: info about a data, but not part of it
- **attribute**: couple name/value
- **superblock**: area where a fs stores its critical info
- **inode**: place to store the metadata of a file
- **dentry**: holds inode's relation, allows fs traversal

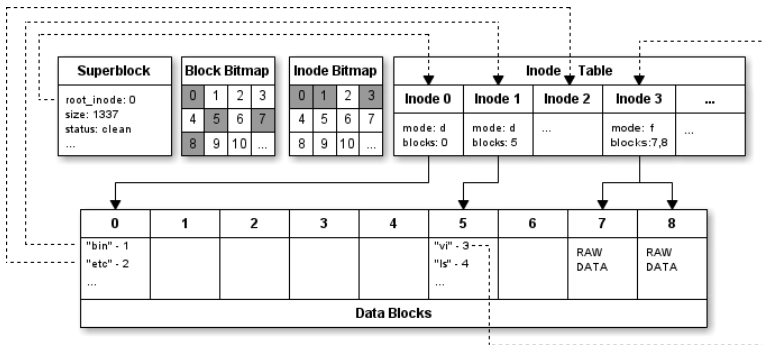
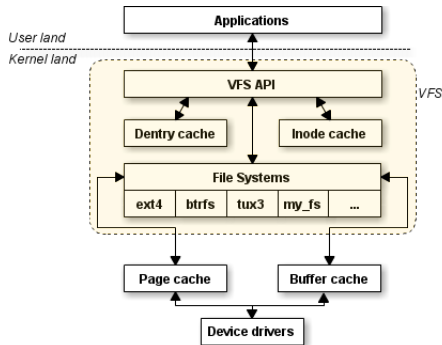
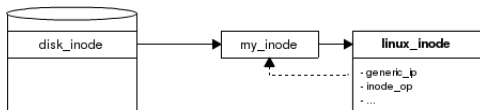


Figure: Relations between superblock, inodes, blocks, ...

- Keep track of available filesystems
- Provide an uniform interface
- Reasonable generic processing for common tasks
- Common I/O cache
  - Page cache
  - I-node cache
  - Buffer cache
  - Directory cache



- Linux defines generic function and structure but doesn't know anything about our fs
- Linux uses composition to store the fs structs
- Each struct contains a pointer to many member functions



Read inode:

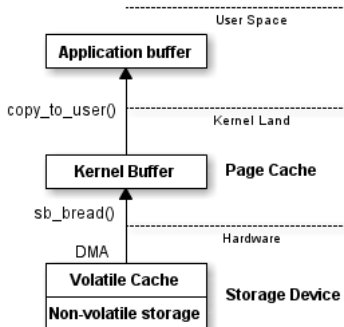
- read disk\_inode from the media
- create a memory representation of it (my\_inode)
- alloc. a linux\_inode and make generic\_ip point to my\_inode
- return the linux\_inode to let the vfs manage it

- Keep disk-backed pages in RAM
- Implemented with the paging memory management
- It uses unused areas of memory.

```
1 42sh> free -m
2          total    used    free   shared  buffers   cached
3 Mem:      2947    2529    417     156     811     709
4 -/+ buffers/cache: 1007 1939
5 Swap:      1953         0    1953
6
```

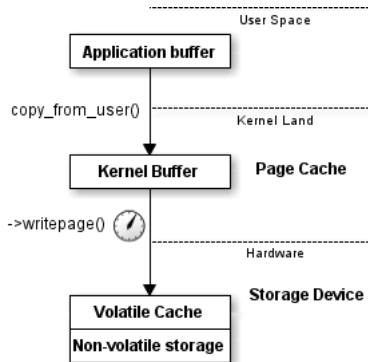
## Reading:

- read syscall
- if in the cache, retrieve it;
- otherwise read it from the device and add it to the cache



## Writing

- Copy buf to the page cache
- Mark the page as dirty
- The kernel periodically transfers all the dirty pages to the device



Why delay the write operations ?

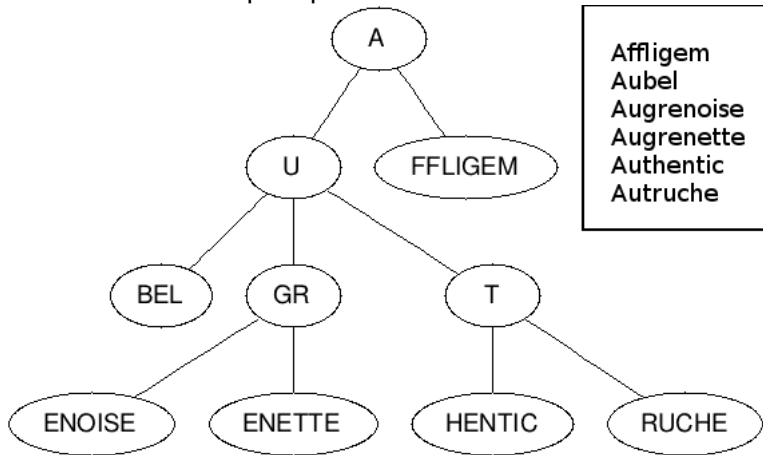
- Temporal locality
- Seek optimization
- Group operations

You can bypass the cache by using **O\_DIRECT**.



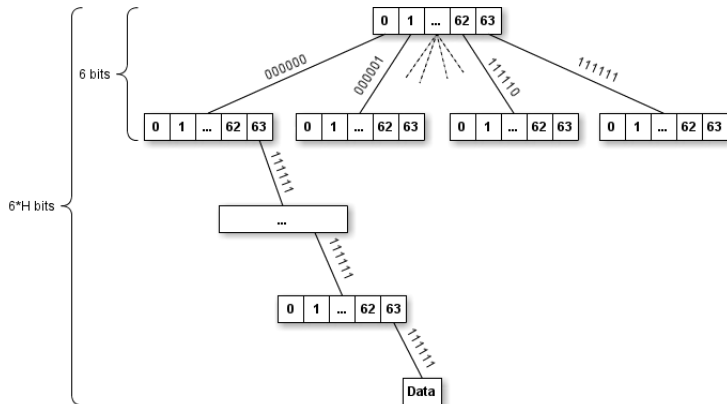
- When free memory shrinks below a specified threshold
- When dirty data grows older than a specific threshold
- Tunable parameters in `/proc/sys/vm/`
- You can change the default I/O scheduler
- If more than a threshold percent of a process's address space is dirty, processes must wait for the I/O scheduler to flush the cache

```
1 cat /proc/sys/vm/dirty_expire_centisecs
2 3000
3 cat /proc/sys/vm/dirty_background_ratio
4 10
5 cat /proc/sys/vm/dirty_ratio
6 40
7
```

**Radix Tree:** a compact prefix tree

## Radix Tree: a compact prefix tree

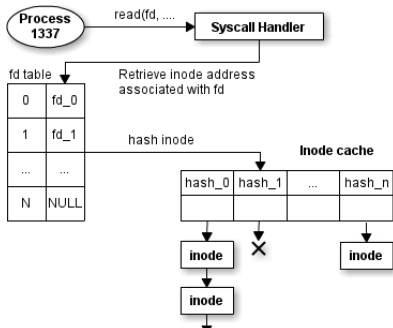
- Wide and shallow
- Each node contain 64 slots
- Each level is a 6 bits prefix



Additional feature: ability to associate tags with specific entries (to mark a page as dirty or under writeback for example) and retrieve them all easily.

## Inode cache:

- Keep recently accessed file i-nodes
- The kernel retrieve the inode from the fd table of the application's address space
- Implemented as an open chain hash table, with blocks linked into a LRU lists
  - Used and dirty
  - Used and clean
  - Unused



**d-cache:** speed up accesses to commonly used directories

- Implemented as an open chained hash table, also linked into a LRU list
- Negative dentry for failed lookups
- Prehash with name, rehash with the dentry parent's address

**Block cache:** interfaces with block devices, and caches recently used meta-data disk blocks. One LRU cache per-CPU.

- The array is sorted, newest buffer is at `bhs[0]`
- Discards the least recently used items first
- Implemented as an array of size 8 (caching 8 pages)

## Section 3

# Logging and Journaling



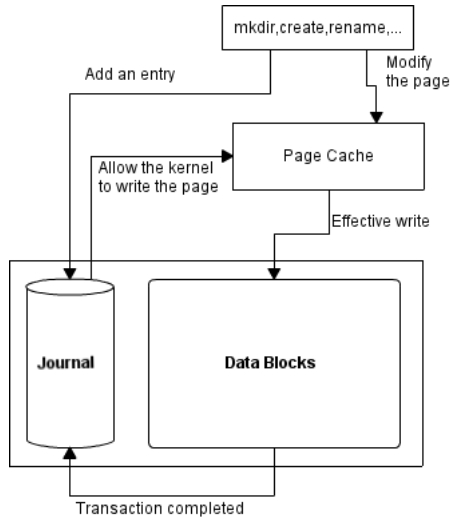
Removing a file :

- Remove its directory entry
- Mark the inode as free
- Mark data blocks as free

A crash between one of these steps leaves the fs in an inconsistent state, and thus needs to be fully checked (fsck)

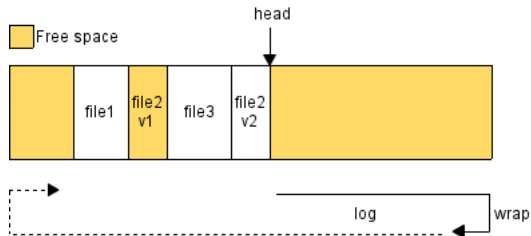
## How to avoid partially written transactions ?

- **Transaction:** complete set of modifications made to the disk during one operation
- **Journal:** Fixed-size contiguous area on the disk (circular buffer)
- Writing to disk:
  - Add an entry to the journal
  - Allow the write to happen on disk
  - Mark the entry as completed
- If an entry is not completed when mounting, replay it



- ⊕ consistency of metadata
- ⊕ faster than fsck
- ⊖ data consistency is not ensured
- ⊖ redundancy of metadata writes

- The whole system data is structured in the form of a circular log
- Avoid writing data twice
- Copy-On-Write, mark the old version as free and write at the end of the log



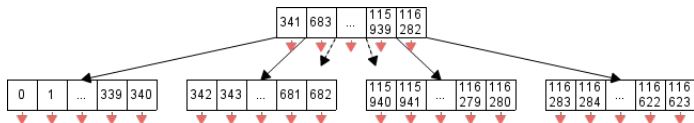
Linux File  
Systems

Adrien  
'schischi'  
Schildknecht

- ⊕ sequential writes
- ⊕ avoid redundant writes
- ⊖ slow random reads

## Section 4

# Real FS design



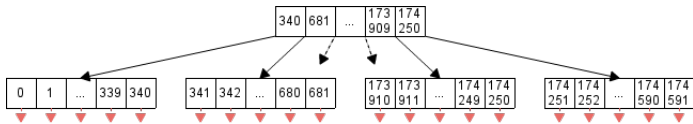
```
1 struct btree_val {
2     int key;
3     void *data;
4 } typedef btree_val;
5 //sizeof(btree_val) = 8
6
```

```
1 struct btree {
2     btree_val values[N];
3     btree *children[N+1];
4 } typedef btree;
5 //sizeof(btree)=8*N+(N+1)*4
6
```

$$4096 \geq N * 8 + (N + 1) * 4$$

$$N = 341$$





```
1 struct bptree {  
2     int key[N];  
3     bptree *children[N+1];  
4 } typedef bptree;  
5 //sizeof(bptree)=4*N+(N+1)*4  
6
```

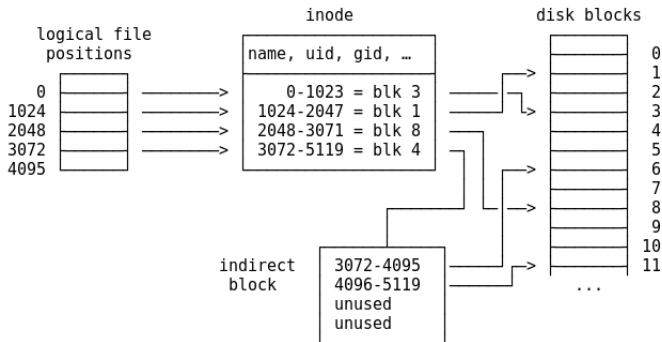
```
1 struct bptree_leaf {  
2     struct {  
3         int key;  
4         void *value;  
5         bptree_leaf *nxt; //opt  
6     } values[M];  
7 } typedef btree_leaf;  
8 //sizeof(btree)=12*N  
9
```

$$4096 \geq 4 * N + (N + 1) * 4$$

$$N = 511$$

$$4096 \geq 12 * M$$

$$M = 341$$



- A chunk of blocks instead of a single block
- Still affected by fragmentation

```
1 struct ext3_extent {
2     __le32 ee_block;    /* first logical block extent covers */
3     __le16 ee_len;     /* number of blocks covered by extent */
4     __le16 ee_start_hi; /* high 16 bits of physical block */
5     __le32 ee_start;   /* low 32 bits of physical block */
6 };
7
```

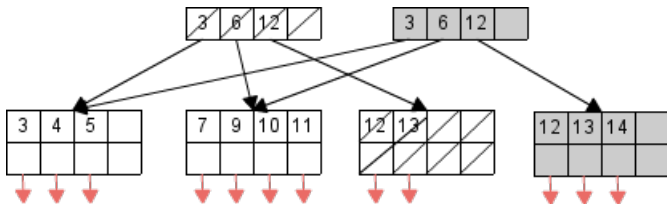
- Inode table
- Extents
- Journal
- Delayed block allocation
- Multi block allocator
- Online defragmentation
- Inline data
- Htree (a variant of B+tree)

Group 0 Padding	ext4 Super Block	Group Descriptors	Reserved GDT Blocks
1024 bytes	1 block	many blocks	many blocks
Data Block Bitmap	inode Bitmap	inode Table	Data Blocks
1 block	1 block	many blocks	many more blocks

- Basically same features as ext4 (extents, inlining, ...)
- Copy On Write metadata and data
- Transparent compression

A B+tree containing a generic key/value pair storage.

- The same btree is used for all metadata



## Section 5

# Conclusion

```
1 #define MEGA(S) ((S) * 1024 * 1024)
2
3 int main(int argc, char *argv[]) {
4     char buf[4096];
5     int fd = open("/home/schischi/foo", O_CREAT | O_WRONLY, 0660);
6
7     if (argc == 2 && !strcmp(argv[1], "-f"))
8         if (fallocate(fd, 0, 0, MEGA(700)) != 0)
9             return 1;
10    for (int i = 0; i < MEGA(700) / sizeof (buf); ++i)
11        write(fd, buf, 4096);
12    write(fd, buf, MEGA(700) % sizeof (buf));
13
14    unlink("/home/schischi/foo");
15    return 0;
16 }
17
```

```
1 $ repeat 100; ./a.out
2     ./a.out  0.01s user 1.46s system 18% cpu 8.018 total
3
4 $ repeat 100; ./a.out -f
5     ./a.out -f  0.00s user 1.01s system 13% cpu 7.440 total
6
```



Linux File  
Systems

Adrien  
'schischi'  
Schildknecht

**Questions ?**

`schischi@lse.epita.fr`

`schischi - irc.rezosup.org`

- FS design

- Book "Practical File System Design" by Dominic Giampaolo

- VFS

- <http://git.kernel.org/cgit/linux/kernel/git/torvalds/linux.git>
- <http://lwn.net/Kernel/Index/>

- Journaling, logging

- <http://pages.cs.wisc.edu/~remzi/OSTEP/file-lfs.pdf>
- <http://research.cs.wisc.edu/wind/Publications/sba-usenix05.pdf>

- Ext4

- [https://ext4.wiki.kernel.org/index.php/Ext4\\_Design](https://ext4.wiki.kernel.org/index.php/Ext4_Design)
- <http://www.ibm.com/developerworks/library/l-anatomy-ext4/>

- Btrfs

- <http://video.linux.com/videos/chris-mason-btrfs-file-system>
- <http://atrey.karlin.mff.cuni.cz/~jack/papers/lk2009-ext4-btrfs.pdf>